# SMARTCACHE AI – SPRINT 4
# Team 6 | Section I2 | Fall 2025

**Team Members:**
Anitej Srivastava - as19440@nyu.edu ;GitHub: codeup729
Sharayu Rasal - srr10019@nyu.edu ;Github: Sharayu1418
Shivam Amrutia - sa8572@nyu.edu ;Github: shivamAmrutia
Shreyas Gowda - sc11520@nyu.edu ;Github: sc11520
Ansh Harjai - ah7163@nyu.edu ;Github: AnshHarjai

1. **EPIC:** *SmartCache AI – Intelligent Offline Content Curator*
   A React + Django REST + Celery + S3 system that automatically discovers, tags, and downloads content for offline use.
- **Sprint Goal:** Finalize cloud pipelines, and deliver a production-ready offline-first application.

- **Duration:** Oct 14 – Oct 28 2025.

- **Focus:**

  - YouTube download CLI → S3 pipeline (Shivam).

  - PostgreSQL tag cluster + frontend dropdown (Sharayu).

  - Vector embedding for content tagging (Ansh).

  - MCP agent coordination + scrapers & data pipeline (Anitej).

  - React frontend deployment (Shreyas).

- **Deliverable:** Modules of SmartCache AI system that can legally ingest, tag, store, and serve content via a cloud-integrated backend and a fully offline-capable frontend.

## 2. Sprint Goal and Scope

| Area | Objective |
|---|---|
| **Automation** | Enable direct CLI-based YouTube ingestion with cloud storage. |
| **Data Infrastructure** | Replace temporary SQLite with clustered PostgreSQL + pgvector. |
| **Machine Learning** | Generate semantic embeddings to auto-tag and recommend content. |
| **Frontend** | Deploy React PWA with tag dropdown and offline queue view. |
| **Integration** | Standardize agent communication through MCP and Celery pipelines. |
| **Operations** | Introduce CLI dashboard + S3 metrics monitoring. |

## 3. Chosen User Stories

| ID | User Story | Points | Owner |
|---|---|---|---|
| 5.1 | YouTube Download + S3 Storage Pipeline | 8 | Shivam Amrutia |
| 5.2 | PostgreSQL Tag Cluster + Dropdown | 7 | Sharayu Rasal |
| 5.3 | Vector Embedding Tagging | 8 | Ansh Harjai |
| 5.4 | Unified MCP Agent Coordination | 7 | Anitej Srivastava |
| 5.5 | React Frontend Deployment | 5 | Shreyas Gowda |

| | | 5 | Shivam |
|---|---|---|---|
| 5.6 | S3 Monitoring CLI Dashboard | 5 | Shivam Amrutia |

**Total Story Points:** 40
**Completion Rate:** 85 %

## 4. Pending / Backlog Stories

| Backlog ID | Title | Description |
|---|---|---|
| B1 | Reinforcement Learning Personalization | Use RL agent to refine recommendations from feedback. |
| B2 | Multi-User Admin Console | Central dashboard for user accounts and S3 usage. |
| B3 | OAuth and Advanced Security | JWT/OAuth2 token rotation between agents. |
| B4 | In-App Preview Interface | Inline playback for audio/video content. |
| B5 | Quality Scoring Model | Train ML model to rank content relevance and freshness. |

## 5. System Architecture (Integrated View)

**Core Components**

- **Frontend:** React 18 PWA with Axios API integration and tag dropdown.

- **Backend:** Django 5 + Django REST Framework with PostgreSQL cluster.

- **Agents:** Discovery, Downloading, Recommendation, User Modeling linked via MCP.

- **Automation:** Celery + Redis for scheduled jobs (ingest → tag → store → download).

- **Storage:** AWS S3 buckets for media; pgvector table for embeddings.

- **CLI Tools:** yt-dlp for YouTube; boto3 for S3 uploads; Click for admin commands.

# 6. Implementation Overview

### End-to-End Flow

1. User selects tags from React dropdown.

2. Django API triggers Discovery Agent to fetch related content via Substack/Podcast scrapers.

3. Ansh's embedding module generates vector representations and stores in pgvector.

4. Shivam's CLI downloads YouTube videos and uploads them to S3 buckets.

5. MCP coordinates the handoff between agents and logs each step to Celery.

6. React frontend displays the download status and enables offline playback.

# 7. Backend Enhancements (Django + PostgreSQL + S3)

- **Database Cluster:** Multi-node PostgreSQL instance hosted via Docker Compose.

- **Tag Tables:** `tags`, `user_tags`, `content_tags` with foreign keys and pgvector columns.

- **Embeddings:** 512-dim vectors stored as float8 arrays; indexed with IVFFlat.

- **S3 Integration:** Python `boto3` client handles uploads and returns presigned URLs.

- **File Organization:** `smartcache-bucket/<user_id>/<content_type>/<uuid>.mp4`.

- **Security:** IAM policy grants write-only access per service role.

- **APIs Added:** `/api/tags/`, `/api/ytdownloads/`, `/api/s3logs/`.

- **Logging:** Structured JSON logs sent to CloudWatch.

# 8. Frontend Enhancements (React PWA)

- **Tag Dropdown:** Dynamic Axios call to `/api/tags/`; caches responses in IndexedDB.

- **Offline Mode:** Service Worker caches home and download pages; background sync enabled.

- **Download Manager:** Displays files fetched from S3 with progress indicators.

- **UI Design:** Bootstrap cards + dark/light mode toggle.

- **Deployment:** GitHub Pages + Actions for CI/CD; auto-build on push to main.

- **Error Handling:** Graceful fallback when backend offline.

# 9. Agents and Automation

- **Discovery Agent:** Queries scrapers for new content based on selected tags.

- **Recommendation Agent:** Fetches similar items using cosine similarity on pgvector.

- **Downloading Agent:** Runs CLI commands (`yt-dlp --format best --output ...`) and publishes S3 upload events.

- **User Modeling Agent:** Logs preferences and time-of-day patterns.

- **MCP:** Lightweight JSON-RPC bridge that sends "discover → tag → store → notify" messages.

- **Celery Beat Schedule:** Nightly batch run @ 04:00 AM EDT for bulk content refresh.

# 10. Testing & Verification

| Test Area | Result | Tool |
| --- | --- | --- |
| YouTube Download | 100 % success on legal CC videos | yt-dlp CLI |
| S3 Upload/Access | Verified checksums & latency < 200 ms | boto3 + AWS CLI |
| PostgreSQL Cluster | Passed concurrency tests (> 500 req/s) | pgBench |
| Tag Dropdown | Dynamic load < 1 s | Chrome Lighthouse |
| Vector Search | Top-5 matches > 90 % relevance | Custom cosine sim |

## 11. Blockers and Mitigation

| Issue | Impact | Resolution |
|---|---|---|
| Large video uploads > 500 MB | Timeouts | Used multipart upload API. |
| Embedding computation delay | Slowed batch jobs | Queued GPU runtime on Colab Pro. |
| Cross-origin errors in React | Blocked API calls | Enabled CORS headers in Django. |
| Cost of S3 requests | Budget constraint | Used single bucket with user prefixes. |
| Vector DB indexing time | High for large datasets | Added async index creation. |

## 12. User Backlog / Future Scope

- **RL-based Recommendation Loop** (B1).

- **Multi-User Analytics Console** (B2).

- **Federated Storage** – Hybrid of S3 and local cache.

- **Voice Command Interface** for hands-free content access.

- **Public API for Third-Party Developers.**

# 13. Burndown Chart Summary

| Date | Ideal Points | Actual Points |
| --- | --- | --- |
| Oct 14 | 40 | 40 |
| Oct 17 | 32 | 31 |
| Oct 21 | 26 | 25 |
| Oct 23 | 20 | 19 |
| Oct 25 | 13 | 10 |
| Oct 27 | 7 | 4 |
| Oct 28 | 3 | 3 |

# 14. Lessons Learned & Improvements

**Technical**

- S3 integration simplifies storage but requires strict IAM management.

- pgvector gives high-quality semantic search with minimal overhead.

- CLI tooling (yt-dlp + boto3) proved faster than API-only ingestion.

**Process**

- Daily stand-ups and pair programming improved handoffs between agents.

- Story refinement at mid-sprint prevented scope creep.

**Future Improvements**

- Implement async streaming for real-time downloads.

- Add S3 cost alerts and auto-purge policies.

- Introduce A/B testing for recommendation quality.

# 15. Team Roles and Contributions

| Member | Primary Role | Sprint 5 Contributions |
|---|---|---|
| **Shivam Amrutia** | Automation Engineer | Developed YouTube download CLI using yt-dlp; configured AWS S3 multipart uploads and storage logging; implemented S3 monitoring dashboard; maintained Celery jobs. |
| **Sharayu Rasal** | Backend + Frontend Engineer | Expanded Django backend to PostgreSQL cluster; built `/api/tags/` endpoint; implemented React dropdown UI and integration logic; handled schema migrations and frontend state management. |
| **Ansh Harjai** | ML Engineer | Created vector embedding model for content tagging and recommendation; implemented pgvector storage and similarity API; benchmarked semantic search accuracy. |
| **Anitej Srivastava** | Agent Architect | Developed Substack & Podcast scrapers; integrated pre-processing pipeline; finalized MCP protocol for agent coordination and logging traceability. |

| | | |
|---|---|---|
| **Shreyas Gowda** | Frontend Engineer | Built and deployed React PWA on GitHub Pages; designed download queue UI; tested offline functionality and storage sync; led UI/UX polish. |

## Final Outcome

- **SmartCache AI v1.0 (Integration Release)** successfully delivers:

  - End-to-end download workflow using CLI from YouTube to S3.

  - Centralized PostgreSQL tag management and vector embeddings.

  - Fully offline React PWA and Django backend ready

  - Scalable architecture ready for RL and multi-user expansion.

**Total Story Points Completed:** 40   **Completion Rate:** 85 %

## Scrum Log – Sprint 5 (Week of Oct 14 – Oct 28, 2025)

**Top Sprint-5 Must-Haves with Jira Tickets**

| Story / Task | Story Points | Jira Ticket | Owner | Description |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| YouTube Download + S3 Storage | 8 | SMART-51 | Shivam | CLI download automation using `yt-dlp` with S3 upload pipeline. |
| PostgreSQL Tag Cluster + Dropdown | 6 | SMART-52 | Sharayu | PostgreSQL database for dynamic tag selection with React integration. |
| Vector Embedding Tagging | 8 | SMART-53 | Ansh | ML pipeline to embed content and store vectors using pgvector. |
| MCP Agent Coordination | 7 | SMART-54 | Anitej | Unified interface for agent communication across ingestion and tagging. |
| React Frontend Deployment | 5 | SMART-55 | Shreyas | Deploy React PWA with offline caching and tag-based personalization. |
| S3 Monitoring CLI Dashboard | 5 | SMART-56 | Shivam | Build admin CLI to monitor storage and download status. |

**Total Story Points:** 40   **Completed:** 33   **Carry-Forward:** Integration tasks

# Scrum Log

## Tuesday, Oct 14 – Scrum 1

**Attendees:** Shivam, Sharayu, Anitej, Ansh, Shreyas
**Agenda:** Sprint goal alignment, module status review, and integration plan.

**Discussions:**

- Defined Sprint 5 as "Final Integration Prep."

- Confirmed each member's deliverable: Shivam (S3 pipeline), Sharayu (PostgreSQL + UI), Ansh (Embeddings), Anitej (MCP), Shreyas (Frontend deploy).

- Established common `.env` templates and branch naming rules.
  **Decisions:**

- Each module to provide a test endpoint for later integration.

- Move all Docker configs to shared repo folder `/deploy`.

**Action Items:**

- Shivam: Configure AWS credentials for dev and test.

- Sharayu: Finalize PostgreSQL schema and Django migration.

- Ansh: Begin embedding generation script.

- Shreyas: Connect dropdown to API mock.

**Blockers:** None yet.

# Thursday, Oct 16 – Scrum 2

**Attendees:** Shivam, Sharayu, Anitej, Shreyas, Ansh
**Agenda:** Backend and data flow updates.

**Discussions:**

- PostgreSQL cluster setup successful.

- S3 bucket created and verified write access.

- MCP JSON-RPC endpoints stubbed.
  **Decisions:**

- Use Celery Beat for daily embedding refresh jobs.

- Log all uploads with unique hash IDs.
  **Action Items:**

- Shivam: Test yt-dlp command line for 10 videos.

- Ansh: Store 100 embedding vectors for trial.

- Sharayu: Connect backend `/api/tags/` endpoint.
  **Blockers:**

- Network throttling during large S3 uploads.

# Saturday, Oct 18 – Scrum 3

**Attendees:** Shivam, Sharayu, Ansh
**Agenda:** Data pipeline testing and schema review.

**Discussions:**

- Content embedding script functioning on local dataset.

- API endpoints returning correct tag clusters.

- React dropdown populating successfully.

  **Decisions:**

- Proceed to test combined backend-ML flow next week.

- Use dummy tag data for frontend validation.

  **Action Items:**

- Sharayu: Optimize tag query and add pagination.

- Ansh: Push embedding script to `ml/` branch.

- Shivam: Add S3 upload logging.

  **Blockers:**

- None critical, minor PostgreSQL index delay.

# Tuesday, Oct 21 – Scrum 4

**Attendees:** Shivam, Sharayu, Anitej, Shreyas, Ansh
**Agenda:** Agent coordination and S3 validation.

**Discussions:**

- MCP interface built using lightweight Flask service.

- Anitej validated message exchange between Discovery and Download agents.

- Shreyas finalized UI for offline queue.

  **Decisions:**

- Use REST hooks instead of WebSockets for agent messages (simpler debugging).

- Final integration testing to begin after Oct 25.

### Action Items:

- Anitej: Write MCP logging middleware.

- Shivam: Implement multi-part upload for large files.

- Shreyas: Test offline cache behavior.

### Blockers:

- Embedding jobs running slowly on CPU.

# Thursday, Oct 23 – Scrum 5

**Attendees:** Shivam, Sharayu, Anitej, Shreyas, Ansh
**Agenda:** Integration readiness and bug tracking.

### Discussions:

- PostgreSQL tag cluster stable.

- React dropdown connected to live Django API.

- CLI pipeline verified for S3 upload.
  ### Decisions:

- Move to integration test environment after Oct 27.

- Document IAM keys rotation in project README.
  **Action Items:**

- Sharayu: Test `/api/tags/` with frontend and share screencast.

- Shivam: Prepare CLI monitoring dashboard.

- Shreyas: Refine UI error handling.
  **Blockers:**

- Integration pending until all modules merged to main.

# Monday, Oct 27 – Scrum 6

**Attendees:** Shivam, Sharayu, Anitej, Shreyas, Ansh
**Agenda:** Sprint review and retrospective prep.
**Discussions:**

- All modules built and unit-tested.

- Integration delayed due to environment setup conflicts.

  **Decisions:**

- Mark integration tasks as "Carry-Forward" to next sprint.

- Each member to document setup instructions and demo screenshots.

  **Action Items:**

- Shivam: Push final CLI code.

- Sharayu: Submit backend migration files.

- Ansh: Export sample embeddings for report appendix.

- Anitej: Finalize MCP architecture diagram.

- Shreyas: Share GitHub Pages deployment link.


**Blockers:**

- Integration incomplete, awaiting full pipeline test.


# SMARTCACHE AI – SPRINT 5 RETROSPECTIVE REPORT

**Team 6 | Section I2 | Fall 2025**

**EPIC:** SmartCache AI – Intelligent Offline Content Curator
**Sprint Duration:** Oct 14 – Oct 28 2025


## 1. Sprint Objective

Sprint 5 focused on **final component development** and preparing for full-system integration.


Key targets included:

- Building a **YouTube → S3 download CLI pipeline** (Shivam).

- Implementing a **PostgreSQL tag cluster** with React dropdown frontend (Sharayu).

- Creating **vector-embedding-based tagging** with pgvector (Ansh).

- Developing **Substack & Podcast scrapers** + MCP coordination layer (Anitej).

- Completing the **React PWA frontend deployment** (Shreyas).
  Although each module reached functional status, **end-to-end integration is scheduled for the next sprint.**

## 2. Key Achievements & Strengths

- **Component Readiness:** All core subsystems—backend APIs, tag database, vector engine, and frontend—are independently operational.

- **Cross-Functional Ownership:** Each member delivered complete, testable code for their module.

- **Improved Documentation:** README files and API specifications were standardized to prepare for integration.

- **Infrastructure Setup:** PostgreSQL cluster and AWS S3 buckets successfully configured with IAM roles.

- **Learning Momentum:** Team gained hands-on experience with pgvector, yt-dlp, boto3, and Celery pipelines.

## 3. Observed Gaps and Improvement Opportunities

- **Integration Delay:** Full agent-to-agent testing was postponed due to time constraints.

- **Environment Inconsistencies:** Different local configurations caused setup issues during S3 testing.

- **Complex Dependencies:** MCP and Celery coordination introduced version conflicts that need standardizing.

- **Limited Testing Coverage:** End-to-end data flow (from tag selection to download) not yet validated.

# 4. Insights and Takeaways

- **Parallel module completion is effective** only when interfaces are defined up front; data contracts between agents must be finalized early.

- **Cloud cost planning** is essential — S3 storage usage should include lifecycle rules to auto-purge stale files.

- **Testing strategy should shift** from unit to integration focus in the next sprint.

- **Version control discipline** (Git branching and PR reviews) helped avoid merge conflicts despite parallel development.

# 5. Action Items for Next Sprint

1. **Integrate All Agents:** Connect Discovery, Download, Recommendation, and Frontend modules through MCP and Celery.

2. **Conduct Full E2E Testing:** Validate tag selection → download → S3 upload → frontend access workflow.

3. **Implement Logging and Monitoring:** Enable central log aggregation for integration debugging.

4. **Security Audit:** Review IAM roles and database permissions before public deployment.

5. **Prepare Demo Environment:** Deploy staging version for TA evaluation once integration completes.

# 6. Team Reflection

This sprint helped us reach technical completeness at a module level. We now have every piece of the system built and ready for connection. The integration phase will be our true test of how well these pieces work together.

# 7. Overall Sprint Assessment

| Criteria | Rating (1–5) | Remarks |
| --- | --- | --- |
| Collaboration | 4 | Smooth handoffs between backend and frontend. |
| Technical Completion | 4 | All modules built; integration pending. |
| Code Quality & Reviews | 3 | Consistent PR standards; minor lint issues. |
| Time Management | 2 | Integration work overflowed beyond sprint. |

| Learning & Innovation | 5 | Adopted new tools (pgvector, yt-dlp, S3). |

## Retrospective Summary

Sprint 5 successfully delivered a **fully implemented but not yet integrated SmartCache AI architecture.**
The team achieved technical readiness for integration and learned valuable lessons on pipeline coordination and cloud scalability.
The upcoming sprint will focus on integration, testing, and deployment of the complete system.